

1 Precision of genetic relationship estimates based on molecular markers

2

3

4 José F. Barbosa-Neto, Carlos M. Hernández, Louise S. O'Donoughue,

5 Mark E. Sorrells\*

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

---

21 J.B. Barbosa-Neto and M.S. Sorrells, Dep. of Plant Breeding and  
22 Biometry, 252 Emerson Hall, Cornell University, Ithaca, NY 14853-1902;  
23 C.M. Hernandez, Biometrics Unit, 324 Warren Hall, Cornell University,  
24 Ithaca, NY 14853-1902; L.S O'Donoughue, Agric. Canada Res. Branch,  
25 Plant Res. Center, Central Exp. Farm, Bldg. 50, Ottawa, Ontario, K1A  
26 OC6, Canada. Paper number 829 of the Plant Breeding Series.

27 \*Corresponding author (mark\_sorrells@qmrelay.mail.cornell.edu).

## ABSTRACT

Genetic progress through selection is directly related to the amount of variability present in the population and the quality of genes contributed by the parents. Molecular markers can be used for estimating genetic relationship between potential parents. A statistical methodology using the size of a  $(1-\alpha)\%$  confidence interval was developed to determine the precision in the estimation of genetic distance between pairs of cultivars. Precision of relationship estimates was affected by type of genetic index used, number of cultivars, and amount of genetic diversity present in the studied group. The size of the  $(1-\alpha)\%$  confidence interval decreased as the number of RFLP fragments increased. Oat and wheat diversity studies were used to illustrate the methodology.

## INTRODUCTION

Cultivar development is based on the selection of superior individuals in segregating populations. Genetic progress through selection is directly related to the amount of variability present in the population and the quality of genes contributed by the parents; as a consequence, the correct choice of parents can maximize the genetic variability and the average performance of a segregating population. Breeders choose parents mainly based on pedigree information and mean performance.

Genetic relationship between cultivars can be used for selection of parents in a breeding program (Frei *et al.*, 1986). Quantitative and qualitative traits have been used to assess genetic diversity between cultivars in oat (Souza and Sorrells, 1991a; Souza and Sorrells, 1991b) and maize (Smith *et al.*, 1990). Molecular markers have been used to predict genetic diversity and combining ability between potential parents. Smith *et al.* (1990) estimated similarities among a group of elite maize (*Zea mays* L.) inbreds and concluded that restriction fragment length polymorphism (RFLP) information could be used to predict combining ability. On the other hand, Melchinger *et al.* (1990) did not detect a correlation between RFLP based distance and combining ability, although the authors recommended RFLP to assign cultivars to different genetic pools. Martin *et al.* (1995) studied the relationship between hybrid performance and parental diversity measured using polymerase chain reaction (PCR) analysis in wheat (*Triticum aestivum* L. em Thell.). They compared seven hard red spring wheats and concluded that PCR analysis using 27 polymorphic primers was not useful for heterosis prediction. Studies on genetic diversity between cultivars using molecular markers have been conducted in several other crops, such as barley (*Hordeum vulgare* L.)

1 (Melchinger *et al.*, 1994), rice (*Oryza sativa* L.) (Wang *et al.*, 1992), and  
2 soybean [*Glycine max* (L.) Merr.] (Keim *et al.*, 1992).

3 Advantages of DNA marker based genetic relationship over other  
4 measures include high polymorphism and detection of allelic variation by  
5 descent and in state. The disadvantages are the cost and time required for  
6 an RFLP study. A practical method is needed to determine the number of  
7 RFLP required to achieve a given level of precision for relationship  
8 between two cultivars. Typically, 25 (Wang *et al.*, 1992) to 640 (Keim *et*  
9 *al.*, 1992) probe-restriction enzyme combinations have been used. Keim *et*  
10 *al.* (1992), working with soybean, used the standard error of the distance  
11 estimation between a pair of cultivars as a criteria to determine the number  
12 of probes required. Tivang *et al.* (1994) used a coefficient of variation  
13 associated with a genetic distance estimation as a measure of precision.  
14 However, their calculations were not generalized for other distance indexes  
15 and plant species.

16 The objective of this study was to develop a statistical approach to  
17 determine the number of probe-restriction enzyme combinations necessary  
18 to estimate genetic distance between cultivars with a given confidence  
19 interval. The methodology was applied to a sample of 34 wheat cultivars  
20 using RFLP and coefficient of parentage to assess genetic relationship.

## 21 THEORY

### 22 Measures of genetic distance

23 Several indices have been proposed to measure genetic distance. Some  
24 indices, such as Nei's distance, include explicit statements about underlying  
25 genetic models. On the other hand, geometric distances do not involve any  
26 genetic concept (Weir, 1990). They are based on Euclidean distance as  
27 follows:

1 
$$d(\text{cultivar 1, cultivar 2}) = [\sum_i^N (p_{i1} - p_{i2})^2]^{1/2},$$
  
2 where  $p_{i1}$  = allelic frequency of the  $i^{\text{th}}$  RFLP fragment in cultivar 1,  $p_{i2}$  =  
3 allelic frequency of the  $i^{\text{th}}$  RFLP fragment in cultivar 2, and  $N$  is the  
4 number of RFLP fragments scored in both cultivars.

5 In this study, a Euclidean distance index was used in order to  
6 differentiate cultivars; because phylogenetic relationship is not relevant for  
7 cultivars and there was no need for genetic concepts, such as mutation rate  
8 or time of divergence. The distance index (DI) used was:

9 
$$DI = [\sum_i^N (p_{i1} - p_{i2})^2 / N]^{1/2}.$$

10 This index ranges from zero, when cultivars are identical at probed  
11 locations, to one, for completely dissimilar cultivars at these locations. It is  
12 similar to Modified Roger's distance (Wright, 1978); however, RFLP  
13 fragments instead of loci were used due to the phenomenon of multiple  
14 dose of fragments in polyploid species. Multiple fragments hinder  
15 determination of allelic associations because a single probe may hybridize  
16 to loci on more than one chromosome (Sorrells, 1992).

#### 17 Statistical distribution of distance index

18 The real frequency of mismatches ( $\phi$ ) between two cultivars can be  
19 determined only at the level of nucleotide sequencing, where total DNA is  
20 sequenced for both cultivars. RFLP fragments and DI are used to estimate  
21  $\phi$ . For inbred cultivars, the allelic frequency of RFLP fragments for any  
22 cultivar can be either zero (absence) or one (presence). As a consequence,  
23 the value of the expression  $(p_{i1} - p_{i2})^2$  will be either zero or one and assuming  
24 independence between RFLP fragments the summation of  $(p_{i1} - p_{i2})^2$  over  
25 RFLP fragments is distributed as a Binomial ( $N, \phi$ ). The estimator DI is  
26 the sample mean of this binomial distribution which has mean  $\phi$  and  
27 variance  $\phi(1-\phi)/N$ . In addition, by the Central Limit Theorem DI is

1 asymptotically normally distributed with mean  $\phi$  and variance  $\phi(1-\phi)/N$ .  
2 An important assumption that was made to establish the statistical  
3 distribution is that each RFLP fragment were independent and identically  
4 distributed. This is not always true, since an RFLP probe that hybridizes  
5 to a fragment that contains an internal restriction site for the restriction  
6 enzyme used could result in two visible fragments.

7 Sample size required to estimate distance between cultivars

8 In order to determine how many probes have to be surveyed to  
9 estimate distance between two cultivars, it is necessary to establish the  
10 degree of precision required for the genetic distance index (DI). One  
11 approach is to use the size of a  $(1-\alpha)\%$  confidence interval (CI) for DI,  
12 where  $(1-\alpha)$  is the probability that the CI covers the real value of DI.  
13 Since DI is a square root of the observed frequency of mismatches ( $f$ ) and  
14 this transformation function is monotone in the interval  $[0, 1]$ ; it is possible  
15 to estimate a confidence interval for  $\phi$  and transform it for DI (Mood et  
16 al., 1974). Let the size of a  $(1-\alpha)\%$  CI for  $\phi$  be  $2E$ . Since  
17

$$E = (Z_{\alpha/2} \sigma) / N^{1/2},$$

18 where  $E$  = half size of a  $(1-\alpha)\%$  CI for  $\phi$ ,  $Z_{\alpha/2}$  = value of the standard  
19 Normal with  $\alpha/2$  probability,  $\sigma$  = standard deviation of  $f$ , estimator of  $\phi$ ,  
20 and  $N$  = total number of RFLP fragments scored; it follows that after  
21 substituting  $\sigma$  for the asymptotic standard deviation:

$$E = Z_{\alpha/2} [f(1-f)/N]^{1/2},$$

23 where  $f$  is the observed frequency of mismatches. Solving for  $N$  in the  
24 above expression:

$$N = [(Z_{\alpha/2})^2 f(1-f)] / (E^2),$$

26 it is possible to calculate the number of RFLP fragments that need to be  
27 scored in order to estimate the genetic distance between two cultivars,

1 given the desired accuracy (E), the level of certainty ( $\alpha$ ), and the number  
2 of RFLP fragment mismatches observed. Table 1 contains the conversion  
3 values of E for f that correspond to E = 0.10, 0.05, and 0.01 for DI.

4 As an example, assume a comparison between two inbreds. Based on  
5 pedigree, previous research, or a preliminary survey, f is set as 0.25. The  
6 experiment requires a 90% CI for DI with a half size E = 0.10. Using  
7 Table 1 for DI = 0.50 , an estimate for the number of scored fragments  
8 required is:

$$9 \quad N = [1.64^2 \times 0.25 \times (1 - 0.25)] / 0.0980^2$$

$$10 \quad N = 52.51$$

11 Thus, for two cultivars, approximately 53 RFLP fragments must to be  
12 scored to achieve the desired precision.

## 13 MATERIAL AND METHODS

### 14 Germplasm

15 Thirty-four wheat cultivars were surveyed for pedigree and RFLP.  
16 This group consisted of winter and spring types and both public and  
17 proprietary inbreds. In addition to the wheat cultivars, RFLP surveys  
18 from 84 oat (*Avena sativa* L.) cultivars were used in order to exemplify  
19 the statistical methodology proposed. Methodology and results for the oat  
20 study were reported by O'Donoghue et al. (1994).

### 21 Coefficient of parentage

22 Pedigrees were obtained from Zeven and Zeven-Hissink (1976) and  
23 Matthews and Anderson (1995). The pedigree for a cultivar was traced  
24 back to landraces or to parental lines with unknown parentage.  
25 Coefficients of parentage (COP) were estimated as described by Cox et al.  
26 (1985). The following assumptions were made: 1) each parent contributed  
27 equally to the offspring and 2) there was no selection during inbreeding.

1 The relationship between a cultivar and its parents was assumed to be  $r_p =$   
2 0.5 and a re-selection of a previous cultivar or landrace had an assigned  
3 relationship of  $r_p = 0.75$ . Cultivars with unknown origin were assumed to  
4 be unrelated ( $r_p = 0$ ) and the relationship of a cultivar with itself was  $r_p =$   
5 1.

#### 6 RFLP analysis

7 A bulk sample of leaf tissue was harvested for DNA extraction from at  
8 least 10 greenhouse grown plants that were 3-4 weeks old. The tissue was  
9 chilled in liquid nitrogen and ground to fine powder. DNA extraction,  
10 DNA digestion, and Southern blotting were done according to Anderson et  
11 al. (1992). The restriction enzyme Eco-RI was used for DNA digestion.  
12 Hybridization procedures were described by Heun et al. (1991), except  
13 exposures of autoradiograms were 2-12 days.

14 Thirty-four probes were selected from the wheat arm map published  
15 by Anderson et al. (1992) to represent loci distributed over the entire  
16 genome. A barley cDNA library (22 probes), an oat cDNA library (9  
17 probes), and a wheat genomic library (3 probes) were used on the 34 wheat  
18 cultivars. All fragments hybridizing to a probe were scored presence (1)  
19 or absence (0). If all fragments from a probe were monomorphic for the  
20 set of cultivars they were excluded from the study.

#### 21 Statistical analysis

22 Genetic distances for each pair of cultivars were calculated from a  
23 matrix constructed with the RFLP data, where rows contained the  
24 genotypes and columns, the fragments scored. Distance index estimations  
25 were calculated using NTSYS-pc (Rohlf, 1992). The observed frequency  
26 of mismatches ( $f$ ) for both the wheat and oat studies were estimated by  
27 averaging the  $f$  obtained for each pair of cultivars. In addition to the total



group, subgroups with 10 or 20 randomly chosen cultivars were used to provide examples for the proposed methodology.

In the wheat cultivars study, cluster analysis was performed for both COP and DI. Dendograms were constructed based on the average linkage process. DI and COP matrixes were correlated and compared using the normalized Mantel statistic, which assumes bivariate normal distribution, but not independence between pairs of observations.

## RESULTS AND DISCUSSION

### Methodology for RFLP sample size estimation

Confidence intervals were calculated for a range of DI and population sizes (Table 2). For illustration, 90% CI was used because other  $(1-\alpha)\%$  CI can be estimated by changing the  $Z_{\alpha/2}$  value. The number of fragments that need to be scored in an RFLP study increases rapidly as the CI size (2E) decreases. The use of an  $E = 0.01$  requires data for a large number of fragments, even when the number of comparisons is very small. On the other hand, the number of fragments required for an  $E = 0.10$  are feasible for the entire range of DI. The CI size where  $E = 0.05$  has an intermediate position, requiring reasonable sample sizes when comparing two or three cultivars or when the estimated value of DI is over 0.60.

The parameter  $\phi$  reflects the genetic diversity of the population under study. Populations composed of related inbred lines generally have larger DI variances compared to populations of distantly related lines. This is due to the square root transformation on the binomial distribution, which makes the variance of DI inversely proportional to the frequency of mismatches observed ( $f$ ). Fig. 1 depicts the effect of different  $\phi$  on the size of a 90% CI for DI. When comparing lines that are related, many RFLP fragments are monomorphic but results in a larger variance estimate. As a

1 consequence, probes with high polymorphism information content (PIC)  
2 are preferred, because the frequency of mismatches will be higher,  
3 resulting in a lower variance. This variance reflects the estimation error  
4 that would be detected in the case of repetitive sampling with different  
5 RFLP probes.

6 The  $(1-\alpha)\%$  CI estimation is used to assess the degree of precision  
7 attained in the estimation of the real DI between cultivars. When more  
8 than two cultivar are studied, it is important to decide if pairwise or  
9 simultaneous comparison will be used. Pairwise comparison is generally  
10 used in studies that attempt to correlate genetic distances with other  
11 measures, such as coefficient of parentage or heterosis. The correlation  
12 coefficient precision is affected by the number of pairs that are being  
13 correlated and not by the precision of the genetic distance estimator. On  
14 the other hand, when a phylogenetic tree is constructed, precise estimation  
15 of genetic distance between all pair of cultivars studied is a primary  
16 requirement. For simultaneous comparison the number of cultivars affects  
17 the number of fragments necessary for accurately estimating genetic  
18 relationship (Table 2). To extend the results to more than two cultivars  
19 one must consider  $M = [L(L-1)/2]$  different pairwise comparisons between  
20  $L$  inbreds. An optimal number of fragments to achieve a desired precision  
21  $E$  will affect the probability of not achieving this precision for all  $M$   
22 comparisons. Therefore, the  $\alpha$  becomes:

$$1-(1-\alpha')^M.$$

24 The pairwise  $\alpha'$ -level must be such that the expression above is equal to the  
25 desired overall  $\alpha$ -level. For example, when comparing line A against  
26 lines B and C, two comparison are made. If an  $\alpha'$ -level of 0.05 is used in  
27 each comparison, the overall  $\alpha$ -level will be approximately 0.10.

## Estimation of $\phi$

For a set of lines not previously studied, the parameter  $\phi$  is unknown for a population and has to be estimated in order to calculate the number of RFLP fragments that need to be surveyed. There are at least three different ways to obtain an estimate for  $\phi$ . The coefficient of parentage between two inbreds could be used as a preliminary estimate or one might obtain a preliminary estimate from previous studies. A more refined approach would be to survey an initial set of probes on the population under study and to use the observed number of mismatches ( $f$ ) as an estimate for  $\phi$ . This preliminary value could be used to estimate the required sample size and then updated as soon as more information is available for the population.

When comparing two cultivars there is only one possible value for  $f$ ; however, in a study involving more than two cultivars there is one value for each pair being compared. The use of an  $f$  value averaged over all comparisons for the formula to obtain  $N$  seems to be a logical approach. On the other hand, if the researcher wants to be more or less stringent, it is possible to use the minimum or maximum observed value of  $f$ , respectively.

## Generalization to other distance measures

A generalization of this methodology to other distance measures can be done using the Central Limit Theorem to approximate the distribution for the index and a  $(1-\alpha)\%$  CI can be estimated. For of inbred lines, where the allelic frequency is always considered to be one or zero, distance or similarity indexes are based on Bernoulli distribution. On the other hand, when the study is conducted with varieties that are not homozygous, the

1 Bernoulli distribution may not apply and the basic statistical distribution  
2 has to be derived.

3 Depending on the dissimilarity within a population, the number of  
4 RFLP fragments that need to be scored to give the same  $(1-\alpha)\%$  CI for  
5 different distance measures may also be different. For example, the  
6 genetic distance indexes DI and Nei's are non-linear transformations of the  
7 frequency of mismatches parameter,  $\phi$  (Fig. 2). In a population with high  
8 dissimilarity between lines, Nei's distance value tends to be infinite, which  
9 would require a large number of RFLP fragments to achieve a reasonable  
10 confidence interval. On the other hand, the distance index used in this  
11 study, DI, requires larger sample sizes for populations with higher  
12 similarity. Roger's distance and Modified Roger's distance exhibit the  
13 same behavior as DI, because these indexes are based on the same  
14 Euclidean distance but are weighted by the number of loci scored rather  
15 than the number of RFLP fragments.

#### 16 Uses of the proposed methodology

17 The proposed methodology allows a researcher to choose an optimal  
18 number of probes to achieve a given level of precision for the estimated  
19 distance index. In order to use the methodology, the average number of  
20 mismatches,  $f$ , has to be calculated as discussed earlier, and an estimation of  
21 the average number of fragments hybridizing to each probe has to be  
22 obtained. The number of fragments per probe does not vary considerably  
23 from collection to collection within the same species; as a consequence, this  
24 value can be obtained from the literature or can be calculated from the  
25 preliminary survey on the collection under study. After attaining the  
26 preliminary sample size, the updated results are used to recalculate the

1 sample size required. This process may be terminated when the  
2 recalculated sample size is similar to the previous one.

3 This statistical methodology does not take into account distribution of  
4 the probes over the genome. In genetic diversity studies, probes should be  
5 chosen to uniformly represent the whole genome, avoiding biases due to  
6 sampling. In addition, only one restriction enzyme should be used in each  
7 study, since polymorphisms obtained with more than one restriction  
8 enzyme may not be independent (Miller and Tanksley, 1990).

#### 9 Oat study: an example

10 O'Donoghue et al. (1994) analyzed genetic diversity in a set of 84 oat  
11 cultivars using RFLP data. The data set containing scores for 358 RFLP  
12 fragments was re-analyzed using DI. A 90% CI and the estimated  
13 minimum RFLP fragment sample size to obtain an  $E = 0.10$  are presented  
14 in Table 3. Subsets of cultivars were selected at random and the average  
15 frequency of mismatches was calculated in order to illustrate the effect of  
16 number of comparisons made in the calculation of the CI. The 90% CI  
17 ranged from  $E = 0.07$  to  $E = 0.10$  when simultaneous comparison was  
18 used. For pairwise comparison the 90% CI was  $E = 0.04$ . The number of  
19 fragments required for an  $E = 0.10$  in the complete set was 60 and 359  
20 RFLP fragments for pairwise and simultaneous comparison, respectively.

#### 21 Genetic relationship for a diverse set of wheat cultivars

22 Table 4 contains characteristics of the 34 wheat cultivars evaluated for  
23 genetic relationship in this study. The hard red grain type was the group  
24 most represented (79%); while soft red and soft white types represented  
25 15% and 6% of the genotypes sampled, respectively. About half of the  
26 cultivars studied were spring types (47%). The number of RFLP  
27 fragments detected per probe ranged from three to 10 with an average of

6.13, and the number of different phenotypes per polymorphic clone from two to eight. Only three out of 34 probes (9%) were monomorphic for all cultivars; 68% of the 190 RFLP fragments scored were polymorphic.

The average genetic distance index within the group of 34 wheat cultivars studied was 0.45. A dendrogram for this set of 34 wheat cultivars using RFLP data suggested two major groups (Fig. 3). The first group was composed of cultivars B17717 through Tascosa and was more variable than the second one, from Bezostaya through Lar193. The cultivar Arthur was equidistant from both groups. Winter and spring types did not cluster in different groups as observed for barley (Melchinger et al., 1994) and oat (O'Donoghue et al., 1994); the intensive use of winter germplasm in the improvement of spring types in wheat may have contributed to this result. All winter types clustered in the first group were related to Turkey (hard red winter cultivar). In the second group some winter types were related to the landrace Mediterranean. Associations between RFLP clusters and other morphological traits, such as grain color and type and plant height, were not observed. Martin et al. (1995) reported an average similarity index (Nei and Li, 1979) of 0.13 in a group of seven hard red spring wheats. The square root of  $(1 - 0.87)$  is equivalent to a DI of 0.36 in their study. The average of 0.45, found in this study, was similar to the oat (DI = 0.40) (O'Donoghue et al., 1994) and barley studies (DI = 0.46) (Melchinger et al., 1994).

Similarly, cluster analysis of the coefficient of parentage (COP) did not separate winter and spring types (Fig. 4). The average COP for the whole group was 0.10, the values ranged from zero to 0.70. It was interesting that COP did not separate winter and spring types to different clusters as observed for oat (Souza and Sorrells, 1991b) and barley (Melchinger et al.,

1994); this result agreed with the RFLP results. Genotypes tended to be grouped according to grain type. Hard red cultivars grouped with at least 12% similarity; however, six hard red genotypes did not share this cluster (Bezostaya, Lac732, Lac519, Lac739, Lar193, and Lar1138). The soft red cultivars Arthur, Caldwell, Pike, and Hart also grouped together with at least 12% similarity. Moro was the cultivar most distantly related by pedigree to any other in the group of 34. This is probably due to the fact that this variety has a large percentage of *Triticum compactum* L. in its parentage.

Comparison of RFLP based genetic distance and coefficient of parentage

The RFLP distance indicated closer genetic relationship between the cultivars than COP. Selection during the process of inbreeding may explain this result, since parents that contribute superior alleles to segregating populations may be over represented in the progeny.

O'Donoughue et al. (1994) reported similar results for oat.

The coefficient of correlation between RFLP based genetic distance and COP ( $r = -0.35$ ) was not significantly different from zero according to the normalized Mantel statistic. By definition, COP- and RFLP-based genetic distance are different measures of genetic relationship between two or more genotypes. Coefficient of parentage measures genetic similarity by descent and is subject to several sources of error. Assumptions about equal contribution from each parent, absence of selection during inbreeding, and unrelatedness of ancestor varieties may bias the estimated COP. In addition, unreliable pedigree information or outcrossing for some parental lines may complicate this value. On the other hand, RFLP genetic distance measures the diversity between genotypes by direct sampling and it provides an estimation of alleles alike in state (Graner et al., 1994).

1 Parents that are used only to improve specific traits, such as disease  
2 resistance, may cause a strong bias in the COP estimation, but not in the  
3 RFLP-based distance, thus reducing correlation between the two measures.  
4 For example, in the study of 34 wheat cultivars, many contained the cross  
5 Norin 10/Brevor in the pedigree. This cross is a major source of dwarfing  
6 genes in wheat and it has been used specifically with the purpose of  
7 reducing stature of new cultivars (Dalrymple, 1980). Autrique (1993) and  
8 Graner et al. (1994) also reported low correlation between RFLP based  
9 genetic distance and COP in durum wheat (*Triticum durum* L.) and barley,  
10 respectively. On the other hand, in maize several authors have reported  
11 high correlations between these two measures (Smith *et al.*, 1990; Smith  
12 and Smith, 1991; Messmer et al., 1993). Probably, hybrid cultivar  
13 development in maize has forced plant breeders to maintain distinct  
14 heterotic groups, increasing the diversity between them, which could be  
15 detected by either RFLP or COP. In addition, higher frequency of RFLP  
16 and better pedigree information in maize may have contributed to the high  
17 correlations observed.

#### 18 Precision of DI estimates in the wheat study

19 An approach similar to the one used in the oat example was used for  
20 the wheat study to evaluate the precision obtained with 190 RFLP  
21 fragments scored. Table 5 presents the estimated half size E of a 90% CI  
22 for the complete group of cultivars as well for the subsets containing 10  
23 and 20 cultivars. For the complete group of cultivars the desired precision  
24 of  $E = 0.10$  was attained with 190 RFLP fragments only for pairwise  
25 comparison. In order to attain  $E = 0.10$  for simultaneous comparison it  
26 would be necessary to score at least 146 more fragments, or 24 more



1 probes. When smaller groups of cultivars were compared, the precision of  
2  $E = 0.10$  was obtained.

3 For precise genetic distance estimates (below  $E = 0.05$ ), it is necessary  
4 to score a large number of RFLP fragments. According to the values of  
5 DI estimated in this study, in the oat study (O'Donoghue et al., 1994), and  
6 in the barley study (Melchinger et al., 1994), for pairwise comparisons,  
7 more than 200 fragments were necessary to be scored to  
8 attain an  $E < 0.05$ . When the required precision of DI estimates increases,  
9 as in the case of phylogenetic trees, the required sample size increases to at  
10 least 1200 fragments or approximately 196 probes.

11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27

## REFERENCES

- Anderson, J.A., Y. Ogihara, M.E. Sorrells, and S.D. Tanksley. 1992. Development of a chromosomal arm map for wheat based on RFLP markers. *Theor. Appl. Genet.* 83: 1035-1043.
- Autrique, J.E. 1993. Molecular markers applied to wheat improvement. PhD.diss. Cornell University, Ithaca, NY.
- Bonierbale, M.W., R.L. Plaisted, and S.D. Tanksley. 1993. A test of the maximum heterozygosity hypothesis using molecular markers in tetraploid potatoes. *Theor. Appl. Genet.* 86: 481-491.
- Cox, T.S., Y.T. Kiang, M.B. Gorman, and D.M. Rodgers. 1985. Relationships between coefficient of parentage and genetic similarity indices in the soybean. *Crop Sci.* 25: 529-532.
- Dalrymple, D.G. 1980. Development and spread of semi-dwarf varieties of wheat and rice in the United States. USDA. Agricultural Economic Report 455.
- Frei, O.M., C.W. Stuber, and M.M. Goodman. 1986. Use of allozymes as genetic markers for predicting performance in maize single cross hybrids. *Crop Sci.* 26: 37-42.
- Graner, A., W.F. Ludwig, and A.E. Melchinger. 1994. Relationships among European barley germplasm: II. Comparison of RFLP and pedigree data. *Crop Sci.* 34: 1199-1205.
- Heun, M., A.E. Kennedy, J.A. Anderson, N.L.V. Lapitan, M.E. Sorrells, and S.D. Tanksley. 1991. Construction of a restriction fragment length polymorphism map for barley (*Hordeum vulgare*). *Genome* 34: 437-447.

- 1 Keim, P., W. Beavis, J. Schupp, and R. Freestone. 1992. Evaluation of  
2 soybean RFLP marker diversity in adapted germplasm. *Theor. Appl.*  
3 *Genet.* 83: 205-212.
- 4 Martin, J.M., L.E. Talbert, S.P. Lanning, and N.K. Blake. 1995. Hybrid  
5 performance in wheat as related to parental diversity. *Crop Sci.* 35:  
6 104-108.
- 7 Matthews, D.E. and O.D. Anderson. 1995. GrainGenes, the Triticeae  
8 Genome Database. Internet anonymous ftp. Host probe.nalusda.gov,  
9 directory pub/grains.
- 10 Melchinger, A.E., A. Graner, S. Mahendra, and M.M. Messmer. 1994.  
11 Relationships among European barley germplasm: I. Genetic diversity  
12 among winter and spring cultivars revealed by RFLPs. *Crop Sci.* 34:  
13 1191-1199.
- 14 Melchinger, A.E., M. Lee, K.R. Lamkey, and W.L. Woodman. 1990.  
15 Genetic diversity for restriction fragment length polymorphisms:  
16 relation to estimated genetic effects in maize inbreds. *Crop Sci.* 30:  
17 1033-1040.
- 18 Messmer, M.M., A.E. Melchinger, R.G. Herrmann, and J. Boppenmaier.  
19 1993. Relationships among early European maize inbreds: II.  
20 Comparison of pedigree and RFLP data. *Crop Sci.* 33: 944-950.
- 21 Miller, J.C. and S.D. Tanksley. 1990. RFLP analysis of phylogenetic  
22 relationships and genetic variation in the genus *Lycopersicon*. *Theor.*  
23 *Appl. Genet.* 80: 437-448.
- 24 Mood, A.M., F.A. Graybill, and D.C. Boes. 1974. Introduction to the  
25 Theory of statistics. 3<sup>th</sup> ed. McGraw-Hill Publishing Company, New  
26 York, NY.

- 1 Nei, M. and W. Li. 1979. Mathematical model for studying genetic  
2 variation in terms of restriction endonucleases. Proc. Natl. Acad. Sci.  
3 USA 76: 5256-5273.
- 4 O'Donoghue, L.S., E. Souza, S.D. Tanksley, and M.E. Sorrells. 1994.  
5 Relationships among North American oat cultivars based on restriction  
6 fragment length polymorphisms. Crop Sci. 34: 1251-1258.
- 7 Rohlf, F.S. 1992. NTSYS-pc. Numerical taxonomy and multivariate  
8 analysis system. Exeter software, Setauket NY.
- 9 Smith, J.S.C. and O.S. Smith. 1991. Restriction fragment polymorphisms  
10 can differentiate among U.S. maize hybrids. Crop Sci. 31: 893-899.
- 11 Smith, O.S., J.S.C. Smith, S.L. Bowen, R.A. Tenborg, and S.J. Wall.  
12 1990. Similarities among a group of elite maize inbreds as measured by  
13 pedigree, F1 grain yield, heterosis, and RFLPs. Theor. Appl. Genet.  
14 80: 833-840.
- 15 Sorrells, M.E. 1992. Development and application of RFLPs in  
16 polyploids. Crop Sci. 32: 1086-1091.
- 17 Souza, E. and M.E. Sorrells. 1991a. Relationship among 70 North  
18 American oat germplasms: I. Cluster analysis using quantitative  
19 characters. Crop Sci. 31: 599-605.
- 20 Souza, E. and M.E. Sorrells. 1991b. Relationship among 70 North  
21 American oat germplasms: II. Cluster analysis using qualitative  
22 characters. Crop Sci. 31: 605-612.
- 23 Tivang, J.G., J. Nienhuis, and O.S. Smith. 1994. Estimation of sampling  
24 variance of molecular marker data using the bootstrap procedure.  
25 Theor. Appl. Genet. 89: 259-264.
- 26 Wang, Z.Y., G. Second, and S.D. Tanksley. 1992. Polymorphism and  
27 phylogenetic relationships among species in the genus *Oryza* as

1       determined by analysis of nuclear RFLPs. Theor. Appl. Genet. 83:  
2       565-581.

3       Weir, B. 1990. Genetic data analysis. Sinauer Associates, Inc. Publishers,  
4       Sunderland, MA.

5       Wright, S. 1978. Evolution and the genetics of populations. IV.  
6       Variability within and among natural populations. The University of  
7       Chicago Press, Chicago, IL.

8       Zeven, A. and N.C. Zeven-Hissink. 1976. Genealogies of 14000 wheat  
9       varieties. The Netherlands Cereal Centre, Wageningen, CIMMYT,  
10      Mexico.

## ACKNOWLEDGMENTS

We would like to acknowledge Cargill Seeds Incorporated for partial support for this project and USDA Plant Genome National Research Initiative for subcontract #92-GO161-Cornell of USDA NRI Grant No.92-37300-7550. We thank Dr. Charles McCulloch for valuable suggestions on the statistical methodology.

1

2 Table 1: Half-size of a  $(1-\alpha)\%$  CI for  $\phi$  required to obtain a  $(1-\alpha)\%$  CI for  
 3 DI with half-size E.

DI	f	E		
		0.10	0.05	0.01
0.10	0.01	0.0100 *	0.0087	0.0020
0.15	0.02	0.0223	0.0142	0.0030
0.20	0.04	0.0347	0.0194	0.0040
0.25	0.06	0.0458	0.0245	0.0050
0.30	0.09	0.0566	0.0296	0.0060
0.35	0.12	0.0670	0.0346	0.0070
0.40	0.16	0.0774	0.0397	0.0080
0.45	0.20	0.0877	0.0447	0.0090
0.50	0.25	0.0980	0.0497	0.0100
0.55	0.30	0.1082	0.0548	0.0110
0.60	0.36	0.1183	0.0598	0.0120
0.65	0.42	0.1285	0.0648	0.0130
0.70	0.49	0.1386	0.0698	0.0140
0.75	0.56	0.1487	0.0748	0.0150
0.80	0.64	0.1587	0.0798	0.0160
0.85	0.72	0.1688	0.0848	0.0170
0.90	0.81	0.1788	0.0898	0.0180
0.95	0.90	0.1890	0.0948	0.0190

4 \* E = 0.07 for DI = 0.10 because DI is not defined for negative values of f.

1 Table 2: Sample size for bands to be scored in order to get a 90%  
2 confidence interval of size 2E given the genetic distance index (DI) and the  
3 number of cultivars to be compared.

DI	E	number of cultivars						
		2	3	5	10	20	30	> 30
0.10	0.07 *	266	441	644	909	1138	1297	1584
	0.05	352	582	851	1201	1503	1714	2093
	0.01	6657	11019	16094	22723	28443	32433	39600
0.20	0.10	86	142	207	293	366	418	510
	0.05	274	454	663	937	1173	1337	1632
	0.01	6455	10685	15606	22034	27581	31451	38400
0.30	0.10	69	114	166	235	294	335	409
	0.05	251	416	608	858	1074	1225	1496
	0.01	6119	10129	14793	20887	26145	29813	36400
0.40	0.10	60	100	146	206	258	294	359
	0.05	229	380	554	783	980	1117	1364
	0.01	5648	9349	13655	19280	24133	27519	33600
0.50	0.10	53	87	127	179	224	256	312
	0.05	204	338	494	697	872	995	1215
	0.01	5043	8348	12192	17214	21548	24571	30000
0.60	0.10	44	73	107	151	189	216	263
	0.05	173	287	419	592	740	844	1031
	0.01	4303	7123	10404	14689	18387	20967	25600
0.70	0.10	35	58	85	119	149	170	208
	0.05	138	228	334	471	589	672	821
	0.01	3429	5676	8291	11706	14652	16708	20400
0.80	0.10	25	41	59	84	105	120	146
	0.05	97	161	235	332	416	474	579
	0.01	2421	4007	5852	8263	10343	11794	14400
0.90	0.10	13	21	31	44	55	63	77
	0.05	51	85	124	175	219	250	305
	0.01	1278	2115	3089	4361	5459	6225	7600

4 \* E = 0.07 for DI = 0.10 because DI is not defined for negative values of f.



- 1 Table 3: Distance index (DI), 90% confidence interval size 2E observed,  
 2 and number of RFLP fragments required to obtain an  $E = 0.10$  in  
 3 randomly selected subsets of oat.

# cultivars in		E		# fragments	
subsets	DI	pairwise	simultaneous	pairwise	simultaneous
10	0.45	0.04	0.07	56	193
10	0.32	0.04	0.08	67	228
10	0.43	0.04	0.07		198
10	0.43	0.04	0.07	58	198
10	0.43	0.04	0.07	58	198
20	0.43	0.04	0.08	58	248
20	0.43	0.04	0.08	58	248
20	0.42	0.04	0.08	59	251
20	0.40	0.04	0.08	60	258
20	0.42	0.04	0.08	59	251
84	0.40	0.04	0.10	60	359

- 4 \* Comparisons between 2 cultivars (pairwise) and between 84 cultivars  
 5 (simultaneous).

1 Table 4: Characteristics of 34 wheat cultivars surveyed for pedigree and  
2 RFLP.

Cultivar	Grain Type	Grain Color	Winter/Spring	Origin
Arthur	Soft	Red	Winter	USA
B17717	Hard	Red	Winter	USA
B617	Hard	Red	Winter	USA
B725S	Hard	Red	Winter	USA
Bezostaya	Hard	Red	Winter	Russia
Blackhull	Hard	Red	Winter	USA
Blueboy	Soft	Red	Winter	USA
Caldwell	Soft	Red	Winter	USA
Centurk 78	Hard	Red	Winter	USA
Era	Hard	Red	Spring	USA
Gaines	Soft	White	Winter	USA
Glenlea	Hard	Red	Spring	Canada
Hart	Soft	Red	Winter	USA
Justin	Hard	Red	Spring	USA
Lac172	Hard	Red	Spring	Argentina
Lac377	Hard	Red	Spring	Argentina
Lac450	Hard	Red	Spring	Argentina
Lac519	Hard	Red	Spring	Argentina
Lac611	Hard	Red	Spring	Argentina
Lac699	Hard	Red	Spring	Argentina
Lac721	Hard	Red	Spring	Argentina
Lac732	Hard	Red	Spring	Argentina
Lac739	Hard	Red	Spring	Argentina
Lancota	Hard	Red	Winter	USA
Lar1138	Hard	Red	Spring	Argentina
Lar193	Hard	Red	Spring	Argentina
Moro	Soft	White	Winter	USA
Olaf	Hard	Red	Spring	USA
Pike	Soft	Red	Winter	USA
R654	Hard	Red	Winter	USA
Selkirk	Hard	Red	Spring	Canada
Sturdy	Hard	Red	Winter	USA
Tascosa	Hard	Red	Winter	USA
Triumph 64	Hard	Red	Winter	USA

1 Table 5: Distance index (DI), 90% confidence interval size 2E observed,  
2 and number of RFLP fragments required to obtain an  $E = 0.10$  in  
3 randomly selected subsets of wheat.

# cultivars in		E*		# fragments*	
subsets	DI	pairwise	simultaneous	pairwise	simultaneous
10	0.43	0.05	0.10	58	198
10	0.46	0.05	0.10	56	190
10	0.46	0.05	0.10	56	190
10	0.45	0.05	0.10	56	193
10	0.46	0.05	0.10	56	190
20	0.45	0.05	0.11	56	241
20	0.44	0.05	0.11	57	244
20	0.45	0.05	0.11	56	241
20	0.44	0.05	0.11	57	244
20	0.44	0.05	0.11	57	244
34	0.45	0.05	0.14	56	336

4 \* Comparisons between 2 cultivars (pairwise) and between 34 cultivars  
5 (simultaneous).

1 Fig. 1: Half-size of a 90% CI (E) according to the number of fragments  
2 scored and three levels of  $\phi$  ( $0.25 = O$ ,  $0.50 = \Delta$ , and  $0.75 = x$ ). A = two  
3 cultivars compared; B = > 30 cultivars compared.

4 Fig. 2: Relationship between the frequency of mismatches ( $\phi$ ) with DI ( $\Delta$ )  
5 and Nei's distance (O).

6 Fig. 3: Dendogram of 34 wheat cultivars compared using RFLP based  
7 distance DI.

8 Fig. 4: Dendogram of 31 wheat cultivars compared using coefficient of  
9 parentage.

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

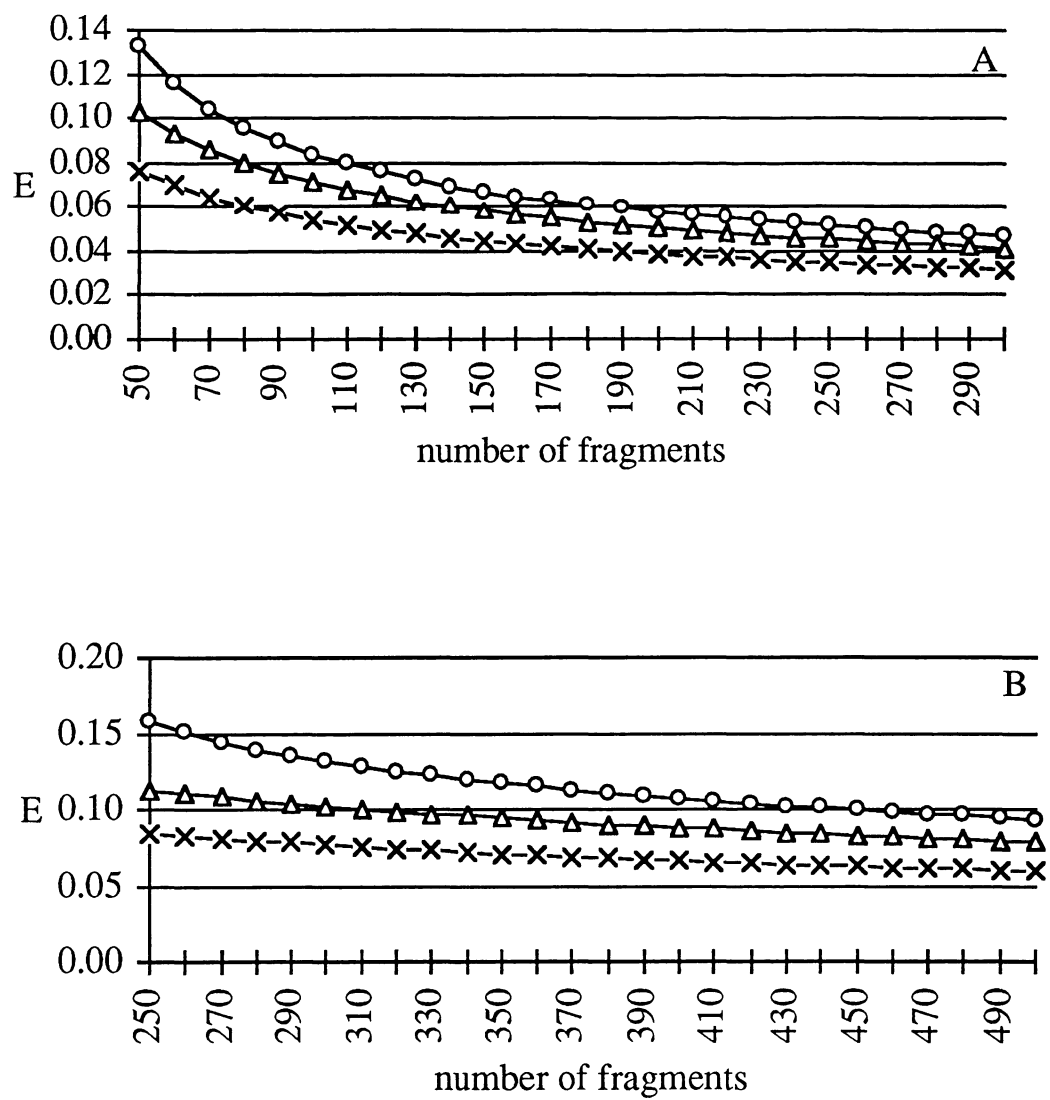


Figure 1

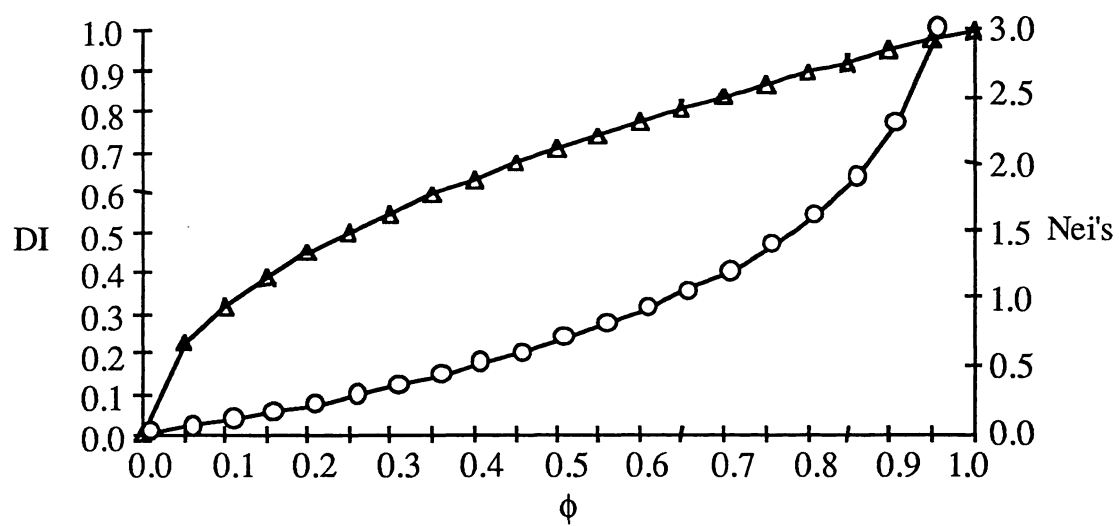


Figure 2

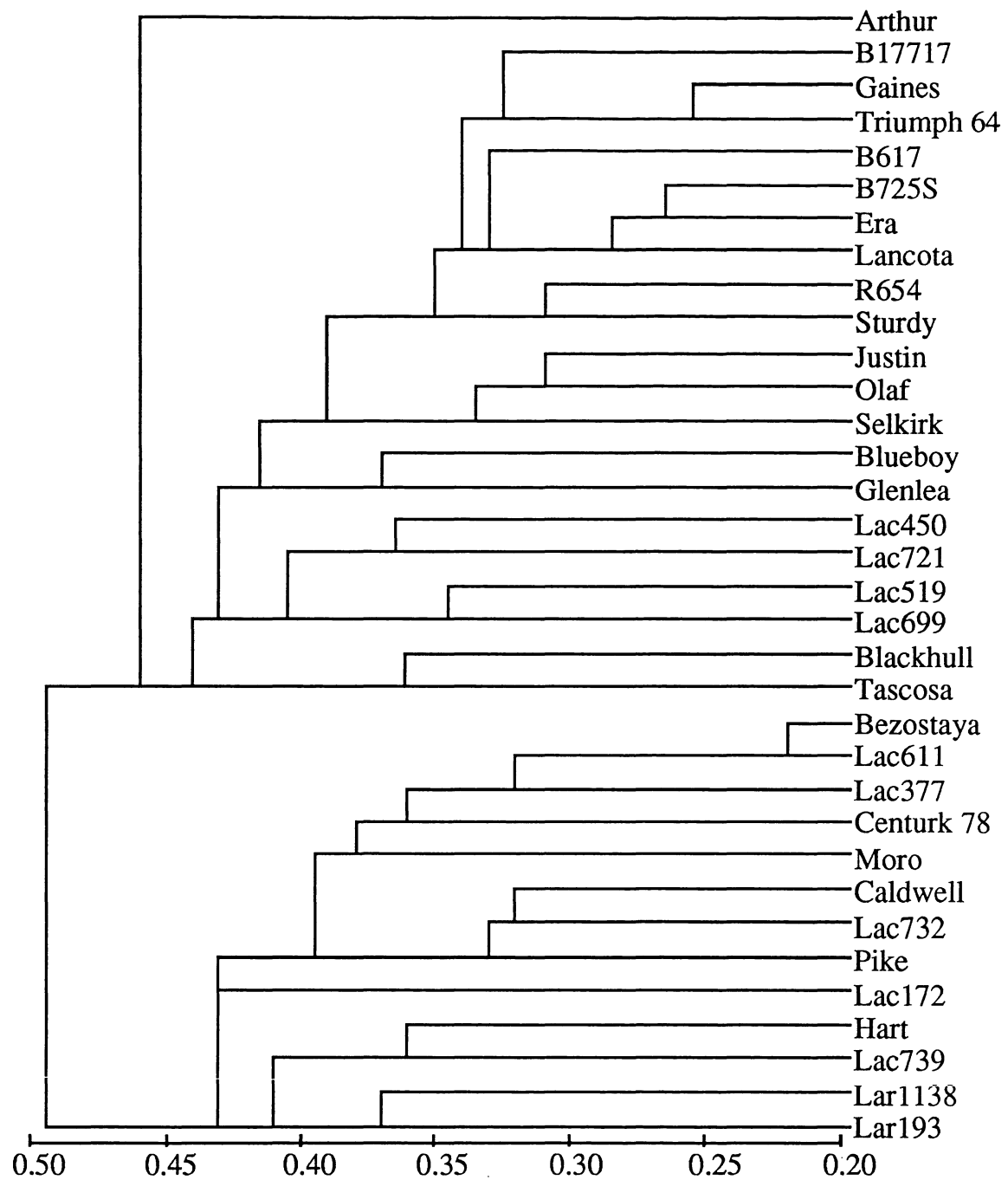


Figure 3

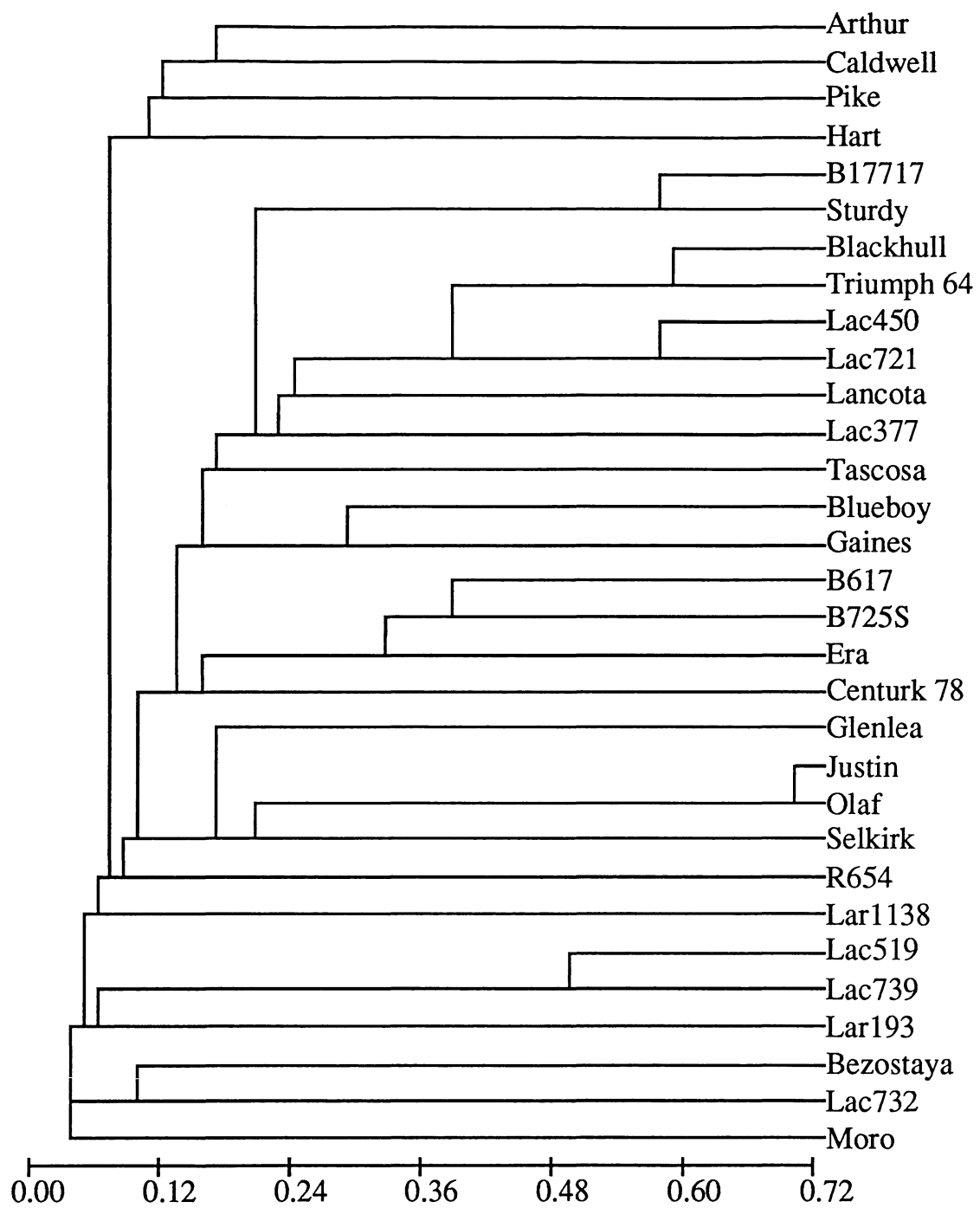


Figure 4